

Fully Convolutional Networks for Handwriting Recognition

Felipe Petroski Such^{*}, Dheeraj Peri^{*}, Frank Brockler[†], Paul Hutkowsky[†], Raymond Ptucha^{*}

^{*}Rochester Institute of Technology, USA, [†]Kodak Alaris, USA

{rwpeec, fps7806, dp1248}@rit.edu

Abstract—Handwritten text recognition is challenging because of the virtually infinite ways a human can write the same message. Our fully convolutional handwriting model takes in a handwriting sample of unknown length and outputs an arbitrary stream of symbols. Our dual stream architecture uses both local and global context and mitigates the need for heavy preprocessing steps such as symbol alignment correction as well as complex post processing steps such as connectionist temporal classification, dictionary matching or language models. Using over 100 unique symbols, our model is agnostic to Latin-based languages, and is shown to be quite competitive with state of the art dictionary based methods on the popular IAM and RIMES datasets. When a dictionary is known, we further allow a probabilistic character error rate to correct errant word blocks. Finally, we introduce an attention based mechanism which can automatically target variants of handwriting, such as slant, stroke width, or noise.

Keywords- Fully Convolutional Neural Networks, Handwriting Recognition, Deep Learning

I. INTRODUCTION

Convolution Neural Networks (CNN) have enabled many recent successes in pattern recognition systems. For example, deep CNNs have become ubiquitous in classification, segmentation and object detection. This success has also been demonstrated in Optical Character Recognition (OCR) systems, where CNNs can predict a sequence of characters from machine generated text with near-perfect accuracy.

Despite a shift towards digitized information exchange, there is still a need for handwritten inputs in many documents such as invoices, taxes, memos, and questionnaires. Intelligent Character Recognition (ICR) is the task of deciphering digitized handwritten text. ICR is quite a bit more challenging than OCR because no two handwritten symbols are identical. ICR handwriting systems can be online or offline. The former records stroke sequences say on a tablet, while the latter has no temporal information. This paper is concerned with offline ICR.

Like OCR, ICR systems first extract lines of text and optionally can segment lines into word blocks separated by white space. ICR systems can feed these word blocks of symbols, often in context with surrounding lines of text, into lexicon-based classifiers. Constraining the output to a lexicon of words can result in very accurate systems, but these same systems unfortunately cannot generalize to commonly occurring sequences such as addresses, phone numbers, and surnames. Removal of dictionary lookup allows the recognition of unbounded dictionaries, but at the expense of decreased accuracy.

Jaderberg et al. [16] proposed a framework which does word level recognition using CNNs on OCR tasks. Poznanski and Wolf [23] used deep CNNs to extract n-grams which feed Canonical Correlation Analysis (CCA) for final word recognition from a fixed vocabulary.

A variant of Recurrent Neural Networks (RNNs), named Long Short Term Memory (LSTM) [14] has shown incredible progress in capturing long term dependencies and have been used in sequence prediction. Some works [21], [5], [22], [9] split an image into segments, then feed the sequence of segments into a RNN to predict a corresponding sequence of output characters. Connectionist Temporal Classification (CTC) [13] further eliminates the need for precise alignment. Xie et al. [31] used CNNs to feed a multi-layer LSTM network for handwritten Chinese character recognition. Similar techniques have also been used for text recognition in natural imagery [25], [26].

A combination of convolutional networks and recurrent neural networks have been used by [28], [30] which use convolution layers for feature extraction and recurrent layers as sequence predictors. [30] performed ICR at the paragraph level to include language context.

Fully Convolutional neural Network (FCN) methods [19], [7] take in arbitrary size images and output pixel level classification. Extracted portions of handwritten text have arbitrary length and can benefit from FCN methods. By using a CNN to estimate the number of symbols in a word block, word blocks can be resized to a canonical representation tuned to a FCN architecture. Knowing the average symbol width, a FCN model can perform accurate symbol prediction without CTC post processing.

This paper proposes a method to obtain character based classification without relying on predefined dictionaries or contextual information. We believe our method is the first that can reliably predict both arbitrary symbols as well as words from a dictionary using a single architecture without any pre- or post-processing. The novel contributions of this paper are: 1) Introduction of a dual stream fully convolutional CNN architecture for accurate symbol prediction; 2) Usage of a probabilistic character error rate that calculates a word probability from a sequence of character probabilities; and 3) Creation of a difficult, but realistic block based dataset derived from the recently released NIST single character dataset [1].

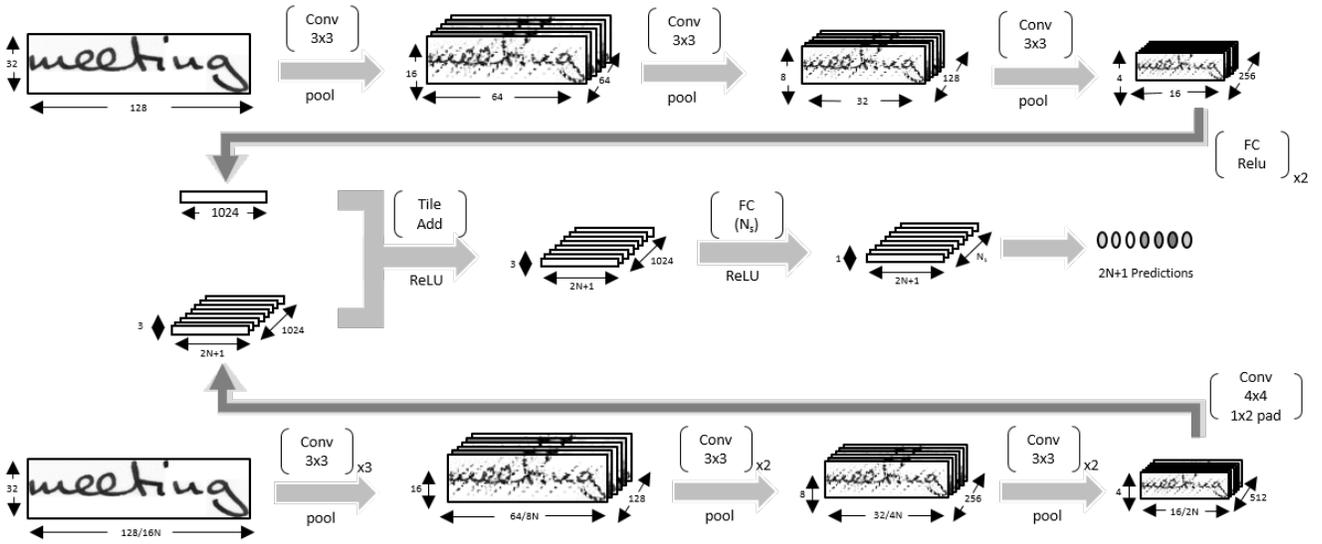


Fig. 1. The Symbol CNN is a fully convolutional model comprised of two networks. The bottom part of the network takes an input image of size $32 \times 128N$, where N is the number of symbols in the word block. The top part of the network processes an image of size 32×128 and adds context during symbol prediction.

II. RELATED WORK

Recent works in scene text recognition [32], [2] use convolutional neural networks and pre-processing techniques like adjusting orientations to predict characters from a scene text. Batuhan et al. [4] proposed an end-to-end framework of recognizing characters from a document by using LSTMs for character segmentation and CNNs for character classification.

One of the advantages of using CNNs is that inputs can be unprocessed data such as raw pixels of an image, rather than extracting specific features [29] or pen stroke grid values. Connectionist Temporal Classification (CTC) removes the need to forcefully align the input stream with character prediction locations [13]. One of the major advantages of the CTC algorithm is that you do not need properly segmented labeled data. Although the CTC algorithm takes care of the alignment of input with the output labels, it can be finicky to train and increases compute complexity during inference.

Huang and Srihari [15] described an approach to separate a line of handwritten text to words. They proposed a gap metrics based approach to perform word segmentation task. Rather than segmenting words, Gader et al. [12] proposed character segmentation utilizing information as you move from background pixels to foreground pixels in horizontal and vertical directions of the character image.

Doetsch et al. [9] proposed hybrid RNN-HMM for English offline handwriting recognition. They introduced a new variant of a LSTM memory cell by using a scalar multiple for every gate in each layer of the RNN. The scaling technique for LSTM gates reduced Character Error Rate (CER) by 0.3%. Bluche et al. [1] compared CNN and traditional feature extraction techniques along with HMM for transcription. Pham et al. [22] proposed Multidimensional RNN using dropout to improve offline handwriting recognition accuracy

by 3%.

Dewan and Srinivasa [8] and Xie et al. [31] used CNNs for offline character recognition of Telugu and Chinese characters respectively. [8] used auto encoders, where the model was trained in a greedy layer wise fashion to learn weights in an unsupervised fashion, then fine-tuned with supervised data.

III. METHODS

An important part of character recognition is Region of Interest (ROI) extraction. Areas of text are extracted from documents using regional based classifiers such as R-CNNs [24] or pixel based segmentation [7]. These regions are split into lines of text using modified XY tree [6]. Each line of text is then split into word blocks, where a word block is a string of symbols separated by white space. The string of symbols could be a word, phone number, surname, acronym, etc. Word block boundaries can again be determined by modified XY tree. Using punctuation detectors [27], it is possible to detect and store a string of word blocks to form sentences, and possible to detect and store a string of sentences to form paragraphs. The FCNs in this research take extracted word blocks as input. Preprocessing such as contrast normalization and deslanting have shown to be effective for handwriting recognition [18] [30], but are not used in this research.

We propose a three staged approach in this paper for the recognition of handwritten characters. In the first stage, we train a CNN to quickly predict the word label for common words such as “the”, “her”, “this”, etc. This model is a typical image classification model where the ground truth classes correspond to discrete words from the training set. We call this a Vocabulary CNN, as it can only predict the common lexicon of words from the training set. The architecture of the Vocabulary CNN is

C(64,3,3)-C(64,3,3)-C(64,3,3)-P(2)-C(128,3,3)-C(128,3,3)-C(256,3,3)-P(2)-C(256,3,3)-C(512,3,3)-C(512,3,3)-P(2)-C(256,4,16)-FC(V)-SoftMax where C(D,H,W) stands for convolution with the dimensions of the filter as $H \times W$ and the depth D . Each convolutional layer is followed by a batch norm and ReLU layer. P(2) represents a 2×2 pooling layer with stride 2. The last layer FC(V) is a fully connected layer with V being the lexicon size. We consider a prediction to be valid if the confidence level is more than 0.7, else we process the block through the next two stages.

In the second stage, we train a CNN to predict the number of symbols in a word block. We refer to this CNN as the Length CNN. This model has a similar architecture to the Vocabulary CNN with additional max-out layers in between convolutional layers. The ground truth classes correspond to the length of the word blocks in the training set. The predicted number of symbols from the Length CNN is used to resize the word block to a canonical representation of $32 \times 128N$, where N is the number of symbols. This resized word block is then passed into the third stage which is a fully convolutional variant of CNN.

In the third stage, a dual stream FCN predicts an arbitrary length sequence of characters from a variable length word block. FCN's are a natural choice for processing variable length word blocks as regular CNNs require fixed length inputs. We call this third stage the Symbol CNN, as it can recognize words at a symbol level. Figure 1 pictorially shows the dual stream architecture used. The top stream takes in an image of size 32×128 which is passed through a series of convolutional layers and a fully connected layer. The fully connected output (fc) size of the top stream is 1024. These features represent the global information in the word block. The second stream consists of only convolutional and pooling layers. The fully connected layer is replaced by a fully convolutional layer which makes $2N + 1$ predictions for a word block of length N . The outputs from both the streams are added and passed through a fully connected layer of size N_s where N_s represents the lexicon of symbols. The final output size is $(2N + 1) \times N_s$. We found that this dual stream model is able to encode rich information about the global and local context of the symbols in the word block.

$2N + 1$ predictions were chosen such that each of the N symbols in a word block are straddled on both sides by a small gap which we refer to as a blank space character. In particular we align the ground truth symbols with blank spaces as shown in Figure 2. In our experiments, we found this to be intuitive and empirically performs well in both dictionary constrained and non-dictionary use cases. There are a total of 123 ground truth symbols which include alphabets (English, French, Spanish, German), numbers, and special characters such as \$, &, period etc.

By not constraining the Symbol CNN output, our model is able to predict any random sequence of characters such as phone numbers, social security numbers, email ID etc. As such, the Symbol CNN is particularly useful to digitize tax, insurance, and medical documents. The dual stream and full convolutional nature of the model enables it to be

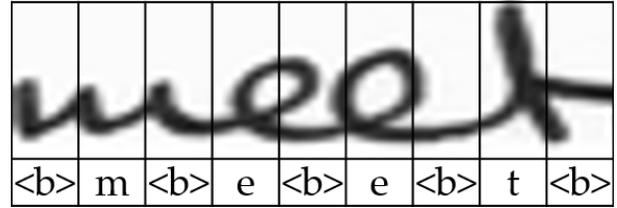


Fig. 2. During training, blank space characters () are added in between ground truth symbols to improve symbol alignment.

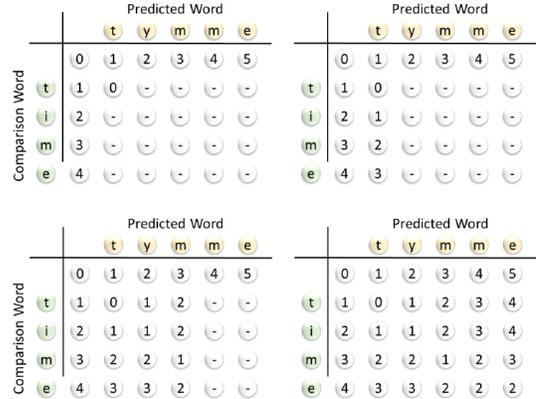


Fig. 3. Calculating CER error between words *tymme* and *time* using dynamic programming. From left to right: after one step, after finishing “t”, after finishing first “m”, and final CER of 2.

quite effective at symbol recognition with or without a word lexicon.

IV. CER AND VOCABULARY MATCHING

We report our results using Character Error Rate (CER):

$$CER = R + D + C \quad (1)$$

where R is number of characters replaced, D is the number of characters deleted and C is the number of correct characters. This can be efficiently computed using dynamic programming and we show an example of the computation in Figure 3. Equation 2 describes the CER computation where $C_{0,0} = 0$ and $CER = C_{l,h}$ where l is the length of the label and h is the length of the prediction. p_i is the i^{th} character of the prediction and L_i is the j^{th} character of the label.

$$C_{i,j} = \min(C_{i-1,j}+1, C_{i,j-1}+1, Diag) \quad (2)$$

where:

$$Diag = \begin{cases} C_{i-1,j-1}, & \text{if } p_i = L_j \\ C_{i-1,j-1} + 1, & \text{otherwise} \end{cases}$$

To improve performance in applications that have a known-limited vocabulary, we applied a CER-based vocabulary matching system using dynamic programming as shown in (3),

$$W(p) = \arg \min_{L \in V} CER(p,L) \quad (3)$$

where V is the vocabulary set and $W(p)$ is the word prediction based on the sequence prediction p . The sequence prediction p is the set of symbols predicted by FCN.

The above method improved CER, but discards most of the information computed with the neural network. An improvement to the above vocabulary matching system, which we refer to as probabilistic CER, uses character probabilities instead of the discrete top character prediction. Equation 4 describes probabilistic CER method.

$$\begin{aligned} C_{i,j} = \min(C_{i-1,j} + 1 - P(p_i = L_j), \\ C_{i,j-1} + 1 - P(p_i = \text{blank}), \\ C_{i-1,j-1} + 1 - P(p_i = L_j),) \end{aligned} \quad (4)$$

Since we now have a method to compute word probabilities from sequence probabilities, we can combine predictions into the same system. We use the frequency of occurrence of a given word $C(L)$ to further improve the vocabulary matching using (5).

$$W(p) = \arg \min_{L \in V} CER(p, L) + \frac{1}{1 + C(L)} \quad (5)$$

V. ATTENTION MODELING

When known variation of handwriting is both common and easily computer generated, data augmentation can be used to generate a family of resized word blocks to predetermined conditions. For example, given an input handwritten word, deslanting algorithms can be used to create a family of inputs, each at a predetermined slant. A similar family of inputs can be created by varying stroke width, noise variation, and paper/background. Given a family of resized word blocks, a vector of attention weights associated with the family of resized word blocks is learned. Using the generated vector attention weight, a single resized word block is formed as a linear combination of resized word blocks, and passed through the CNN in normal fashion.

VI. TRAINING

We used Caffe [17] to perform all experiments. Custom layers were implemented in python to handle ground truth and the blankspace symbol alignment. We trained jointly on IAM [20], RIMES [3] and NIST [1] datasets. All experiments used batch size of 64 (24 IAM, 24 RIMES and 16 NIST samples) and ran for 50k iterations. We used a learning rate of 0.001 to train the network with 0.9 momentum. The learning rate was decreased to 0.0001 after 40k iterations. We also used L2-regularization with $\lambda = 0.0025$.

VII. RESULTS

Results are demonstrated on the IAM, RIMES, and NIST offline handwritten datasets. The IAM [20] dataset contains 115,320 English words, mostly cursive, by 500 authors. This dataset includes training, validation, and test splits, where an author contributing to a training set, cannot occur in the validation or test split. The RIMES [3] dataset contains 60,000 French words, by over 1000 authors. There are

several versions of the RIMES dataset, where each newer release is a super-set of prior releases. We utilize the popular ICDAR2011 release.

The NIST Handprinted Forms and Characters Database, Special Database 19 [1], contains NIST’s entire corpus of training materials for handprinted document and character recognition. Each author filled out one or more pages of the NIST Form-based Handprint Recognition System. It publishes handprinted sample forms from 810,000 character images, by 3,600 authors.

A. IAM Results

We first test our system on the IAM English handwritten dataset. We fine-tuned our model on IAM and it achieves CER of 4.43%. Table I highlights the results of our model with dictionary correction which is quite competitive to the current leaders of this dataset. Table II shows examples of predictions obtained on the IAM dataset using only the FCN model without dictionary correction.

TABLE I

COMPARISON OF RESULTS ON IAM DATASET TO PREVIOUS METHODS.

| Model | WER | CER |
|--------------------------|-------------|-------------|
| Dreuw et al. [10] | 18.8 | 10.1 |
| Boquera et al. [11] | 15.5 | 6.90 |
| Kozielski et al. [18] | 13.30 | 5.10 |
| Bluche et al. [5] | 11.90 | 4.90 |
| Doetsch et al. [9] | 12.20 | 4.70 |
| Our work | 8.71 | 4.43 |
| Voigtlaender et al. [30] | 9.3 | 3.5 |
| Poznanski and Wolf [23] | 6.45 | 3.44 |

Dreuw et al. [10] showed that competitive results can be obtained by hybrid approaches of MLP and Gaussian HMMs. Kozielski et al. [18] used a novel HMM based system. Dreuw et al. [10] and Boquera et al. [11] use a hybrid neural network and Hidden Markov Model HMM approach. Bluche et al. [5] used Gaussian HMMs to initialize neural networks and showed that both deep CNNs and RNNs could produce state of the art results. Doetsch et al. [9] uses a custom LSTM topology along with CTC alignment. [9], [5] used all words in a sentence and paragraph respectively to provide word context. Poznanski and Wolf [23] used deep CNNs to extract n-gram attributes which feed CCA word recognition. [18], [23], [9], [11] use deslanting, training augmentation, and an ensemble of test samples.

Our work uses a first Vocabulary CNN of 1100 common words. The Symbol CNN uses 123 symbols, and we use probabilistic CER correction. Aside from the probabilistic CER correction, no CTC alignment or CCA post correction was applied. Although our competitive results are not ranked the best, our processing path can work at both the symbol and lexicon level, and we include substantially more symbols than prior methods (e.g. [23] can only recognize upper and lower case Latin alphabet).

B. RIMES Results

We use a vocabulary CNN of 800 common words and fine-tuned our model on RIMES dataset. Table III shows

TABLE II

SYMBOL SEQUENCE PREDICTION ON THE IAM DATASET. THE THIRD EXAMPLE HAS A QUESTIONABLE GROUND TRUTH AND A PREDICTION THAT COULD BE CONSIDERED VALID OUT OF CONTEXT.

| Input | Label | Prediction |
|-------|------------|------------|
| | that | that |
| | had | had |
| | Liverpool | livepool |
| | on | oui |
| | mistaken | mistahon |
| | , | , |
| | implements | implement |
| | least | least |
| | mist | mist |
| | interest | interest |

our model obtained a CER of 2.22% using an ensemble of three models, each with dictionary correction. Table IV shows examples of predictions obtained from one of the models without dictionary correction. In general, errors can be attributed to character ambiguity, segmentation artifacts (sample “effet” contains a comma even though it isn’t part of the label). For most of the examples in Table IV, our model made perfect predictions.

TABLE III

COMPARISON OF RESULTS ON RIMES DATASET TO PREVIOUS METHODS.

| Model | WER | CER |
|-------------------------|-------------|-------------|
| Kozielski et al. [18] | 13.70 | 4.60 |
| Doetsch et al. [9] | 12.90 | 4.30 |
| Bluche et al. [5] | 11.80 | 3.70 |
| Our work | 5.68 | 2.22 |
| Poznanski and Wolf [23] | 3.90 | 1.90 |

TABLE IV

SYMBOL SEQUENCE PREDICTION ON THE RIMES DATASET.

| Input | Label | Prediction |
|-------|--------------|--------------|
| | vous | vous |
| | titre | titre |
| | avancé | avance |
| | effet , | effett |
| | désire | diésiire |
| | téléphone | téléhone |
| | relevés | relves |
| | salutations | salutations |
| | l’expression | l’expression |
| | effectuer | effectuer |

TABLE V

SYMBOL SEQUENCE PREDICTION ON NIST GENERATED SAMPLES.

| Input | Label | Prediction |
|-------|-------------------|------------------|
| | 9/10/1966 | 9/10/1966 |
| | (246)344-9702 | (246)344-9702 |
| | \$8643133 | \$8643133 |
| | Spectrometry | Spectrometry |
| | +6091620 | +6091620 |
| | 92.84.8A.b4.AE.15 | 92.84.8A.b4AE.15 |

C. NIST

While there are several class specific handwritten datasets, both at the character and word level, there is no large handwritten dataset that concentrates on word blocks of arbitrary symbols. To test the performance of our model on generic word blocks made of arbitrary symbols, we created a new symbol recognition dataset by stochastically combining the NIST individual character images into realistic word blocks. Images of hand printed text are simulated by extracting character images from a randomly selected writer in the NIST dataset and concatenating them into word blocks of random dictionary words, random strings of alphanumeric characters, or random strings of numeric characters. In addition, the NIST dataset has been supplemented with handwritten punctuation, mathematical symbols and other common special characters such as the dollar sign, email character and the ampersand to facilitate in generating word block images of common form-field inputs.

The images are further augmented by adding random amounts of displacement, stretching and rotation to each symbol to simulate the natural variability in writer’s penmanship. A random amount of skew is then applied to each concatenated image to vary the slant of the word block. Finally, random amounts of noise and blur are added to simulate the effects of image digitization. Table V shows randomly generated samples from NIST characters and corresponding predictions. Our model achieves 92.4% accuracy on a subset of 12,000 word blocks (consisting English, French and special characters) generated from NIST.

Some of the characters in handwritten word blocks are too hard to decipher and require contextual information for human or machine interpretation. In order to provide more contextual information and recognize characters which are ambiguous due to incomplete pen strokes and inconsistencies in handwriting, we utilized both the contextual path introduced at the top of Figure 1 as well as an attention mechanism. After preprocessing an input word block with a series of known variations such as slant and noise, the attention mechanism formulates a single input as a linear combination of these variations. Although the attention mechanism is intuitive, empirical results have yet to show any significant improvement and were not used in the results in Tables I or III.

VIII. CONCLUSION

We introduce a novel offline handwriting recognition algorithm using a fully convolutional network. Unlike lexicon constrained methods, our method can recognize common words as well as infinite symbol blocks such as surnames, phone numbers and acronyms. Our dual stream architecture along with blank space symbol alignment eliminates the need of complex character alignment methods such as CTC in recurrent based methods. Our fully convolutional model enables processing of arbitrary length inputs and utilizes a large symbol set. Despite not using a word lexicon and handling of a large symbol set, our method performs on par with current state of the art methods on English based IAM, French based RIMES, and NIST arbitrary symbol handwritten dataset. Our model is flexible and can be easily extended to other languages.

REFERENCES

- [1] Nist special database 19, 2nd ed. Accessed: 2017.
- [2] Saad Bin Ahmed, Saeeda Naz, Muhammad Imran Razzak, and Rubiyah Yousaf. Deep learning based isolated arabic scene character recognition. In *Arabic Script Analysis and Recognition (ASAR), 2017 1st International Workshop on*, pages 46–51. IEEE, 2017.
- [3] Emmanuel Augustin, Matthieu Carré, Emmanuèle Grosicki, J-M Brodin, Edouard Geoffrois, and Françoise Prêteux. Rimes evaluation campaign for handwritten mail processing. In *International Workshop on Frontiers in Handwriting Recognition*, pages 231–235, 2006.
- [4] Batuhan Balci, Dan Saadati, and Dan Shiferaw. Handwritten text recognition using deep learning. *CS23In: Convolutional Neural Networks for Visual Recognition, Stanford University Project*, 2017.
- [5] Théodore Bluche, Hermann Ney, and Christopher Kermorvant. A comparison of sequence-trained deep neural networks and recurrent neural networks optical modeling for handwriting recognition. In *International Conference on Statistical Language and Speech Processing*, pages 199–210. Springer, 2014.
- [6] Francesca Cesarini, Marco Gori, Simone Marinai, and Giovanni Soda. Structured document segmentation and representation by the modified xy tree. In *Proceedings of the Fifth International Conference on Document Analysis and Recognition*, pages 563–566, 1999.
- [7] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv preprint arXiv:1606.00915*, 2016.
- [8] Sagar Dewan and Srinivasa Chakravarthy. A system for offline character recognition using auto-encoder networks. In *International Conference on Neural Information Processing*, pages 91–99. Springer, 2012.
- [9] Patrick Doetsch, Michal Kozielski, and Hermann Ney. Fast and robust training of recurrent neural networks for offline handwriting recognition. In *14th International Conference on Frontiers in Handwriting Recognition*, pages 279–284, 2014.
- [10] Philippe Dreuw, Patrick Doetsch, Christian Plahl, and Hermann Ney. Hierarchical hybrid mlp/hmm or rather mlp features for a discriminatively trained gaussian hmm: a comparison for offline handwriting recognition. In *18th IEEE International Conference on Image Processing*, pages 3541–3544, 2011.
- [11] Salvador Espana-Boquera, Maria Jose Castro-Bleda, Jorge Gorbe-Moya, and Francisco Zamora-Martinez. Improving offline handwritten text recognition with hybrid hmm/ann models. *IEEE transactions on pattern analysis and machine intelligence*, 33(4):767–779, 2011.
- [12] Paul D Gader, Magdi Mohamed, and Jung-Hsien Chiang. Handwritten word recognition with character and inter-character neural networks. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 27(1):158–164, 1997.
- [13] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*, pages 369–376. ACM, 2006.
- [14] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [15] Chen Huang and Sargur N Srihari. Word segmentation of off-line handwritten documents. In *Document Recognition and Retrieval XV*, volume 6815, page 68150E. International Society for Optics and Photonics, 2008.
- [16] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Synthetic data and artificial neural networks for natural scene text recognition. *arXiv preprint arXiv:1406.2227*, 2014.
- [17] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*, pages 675–678. ACM, 2014.
- [18] Michal Kozielski, Patrick Doetsch, and Hermann Ney. Improvements in rwth’s system for off-line handwriting recognition. In *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*, pages 935–939. IEEE, 2013.
- [19] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- [20] U-V Marti and Horst Bunke. The iam-database: an english sentence database for offline handwriting recognition. *International Journal on Document Analysis and Recognition*, 5(1):39–46, 2002.
- [21] Farès Menasri, Jérôme Louradour, Anne-Laure Bianne-Bernard, and Christopher Kermorvant. The a2ia french handwriting recognition system at the rimes-icdar2011 competition. In *Document Recognition and Retrieval XIX*, volume 8297, page 82970Y. International Society for Optics and Photonics, 2012.
- [22] Vu Pham, Théodore Bluche, Christopher Kermorvant, and Jérôme Louradour. Dropout improves recurrent neural networks for handwriting recognition. In *14th International Conference on Frontiers in Handwriting Recognition*, pages 285–290, 2014.
- [23] Arik Poznanski and Lior Wolf. Cnn-n-gram for handwriting word recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2305–2314, 2016.
- [24] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [25] Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE transactions on pattern analysis and machine intelligence*, 39(11):2298–2304, 2017.
- [26] Baoguang Shi, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. Robust scene text recognition with automatic rectification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4168–4176, 2016.
- [27] Carlos N Silla and Celso AA Kaestner. An analysis of sentence boundary detection systems for english and portuguese documents. In *International Conference on Intelligent Text Processing and Computational Linguistics*, pages 135–141. Springer, 2004.
- [28] Zenghui Sun, Lianwen Jin, Zecheng Xie, Ziyong Feng, and Shuye Zhang. Convolutional multi-directional recurrent network for offline handwritten text recognition. In *15th International Conference on Frontiers in Handwriting Recognition*, pages 240–245, 2016.
- [29] Øivind Due Trier, Anil K Jain, and Torfinn Taxt. Feature extraction methods for character recognition—a survey. *Pattern recognition*, 29(4):641–662, 1996.
- [30] Paul Voigtlaender, Patrick Doetsch, and Hermann Ney. Handwriting recognition with large multidimensional long short-term memory recurrent neural networks. In *15th International Conference on Frontiers in Handwriting Recognition*, pages 228–233, 2016.
- [31] Zecheng Xie, Zenghui Sun, Lianwen Jin, Ziyong Feng, and Shuye Zhang. Fully convolutional recurrent network for handwritten chinese text recognition. In *23rd International Conference on Pattern Recognition*, pages 4011–4016, 2016.
- [32] Usha Yadav, Satya Verma, Deepak Kumar Xaxa, and Chandrakant Mahobiya. A deep learning based character recognition system from multimedia document. In *Power and Advanced Computing Technologies (i-PACT), 2017 Innovations in*, pages 1–7. IEEE, 2017.